## Efficient Autoregressive Audio Modeling via Next-Scale Prediction





Figure 1: Autoregressive modeling of audio. (a) Next-token prediction: sequential token generation in chronological order (left to right), which aligns with the natural temporal structure of audio; (b) Next-scale prediction: multi-scale token maps are autoregressively generated from coarse to fine scales (lower to higher resolutions). Tokens are generated in parallel within each scale, which reduces about 40x the AR prediction iteration.

## Abstract

Audio generation has achieved remarkable progress with the advance of sophisticated generative models, such as diffusion models (DMs) and autoregressive (AR) models. However, due to the naturally significant sequence length of audio, the efficiency of audio generation remains an essential issue to be addressed, especially for AR models that are incorporated in large language models (LLMs). In this paper, we analyze the token length of audio tokenization and propose a novel Scale-level Audio Tokenizer (SAT), with improved residual quantization. Based on SAT, a scale-level Acoustic AutoRegressive (AAR) modeling framework is further proposed, which shifts the next-token AR prediction to next-scale AR prediction, significantly reducing the training cost and inference time. To validate the effectiveness of the proposed approach, we comprehensively analyze design choices and demonstrate the proposed AAR framework achieves a remarkable  $35 \times$  faster inference speed and +1.33 Fréchet Audio Distance (FAD) against baselines on the AudioSet benchmark.

## 1 Introduction

001

011

012

024

Autoregressive (AR) modeling (Achiam et al., 2023; Sun et al., 2024) has been widely used in the generation domain, which typically involves two steps - token quantization (Esser et al., 2021;

Yu et al., 2021) and next-token prediction (Achiam et al., 2023; Touvron et al., 2023). Specifically, the token quantization aims to convert the inputs to a sequence of discrete tokens and the next-token prediction models the conditional distribution of one token based on previous ones. AR approaches have shown significant success in textual modeling, e.g., large language models (LLMs) (Vaswani et al., 2017; Devlin et al., 2018; Touvron et al., 2023; Achiam et al., 2023) and even visual modeling (Dosovitskiy et al., 2020; Chang et al., 2022). However, despite its effectiveness, AR-based audio generation remains under-explored.

Unlike natural language which is discrete and can be easily tokenized into a short series of tokens, audio demonstrates more challenges to be discretized without losing perceptual quality given its long sequence and continuity nature. Previous approaches (Défossez et al., 2022; Yang et al., 2023a; Kumar et al., 2024; Zeghidour et al., 2021) leverage multi-stage residual quantization (RQ) (Lee et al., 2022) to model the raw waveform with different frequencies. However, the multi-stage RQ will significantly increase the token length, leading to difficulty in the subsequent next-token prediction. Another paradigm (Baevski et al., 2020) focuses on the semantics of the waveform and leverages pre-trained models (e.g., Hubert (Hsu et al., 2021)) to cluster the embeddings in the semantic

084 880

094

100

102

103

104

105 106

108

110

space and then quantize the embeddings based on cluster centers. Though semantic embeddings can successfully reconstruct the waveform, the reconstruction quality and generalization capability are bottlenecked by the pre-trained encoder.

In addition, compared to text and images, audio waveform typically has a much longer sequence length due to the high sampling rates, such that about 960000 sequence length in 1 min audio clip with 16kHz. Since AR models predict tokens in a sequential manner, the inference cost is quadratically correlated to the sequence length, making the AR-based audio generation slow and computationally expensive, as illustrated in Fig. 1 (a).

VAR (Tian et al., 2024; Li et al., 2024a,c), a recent variant of AR models shifts from token-wise to scale-wise AR prediction scheme with a multiscale tokenizer, showcasing improved efficiency and scalability in visual domains. However, applying scale-wise prediction to raw audio generation remains challenging due to the high temporal resolution of audio signals, making efficient audio generation difficult with existing methods. Large token length from the tokenizer will burden the AR modeling. Unlike the 2D visual tokenizer that compresses images in spatial (vertical and horizontal) axes, the audio tokenizer typically only compresses along the temporal axis, making it challenging to achieve a high compression rate. To address this issue, we leverage a trade-off between token length and the residual depth. Specifically, multi-scale quantization reduces the token number in each scale, allowing for more residual layers under the same total token constraint, thereby enhancing performance. This observation highlights the potential of multi-scale designs for tokenizer optimization, especially for audio applications.

In this paper, we explore the Scale-level Audio Tokenizer and Multi-Scale Acoustic AutoRegressive Modeling in audio generation, as illustrated in Fig. 1 (b). On the one hand, to shorten the audio token length, we utilize a scale-level audio tokenizer (SAT) which improves the traditional residual quantization with a multi-scale design and compresses the token length according to the scale index. On the other hand, we further shorten the inference step during the autoregressive prediction. Based on the multi-scale audio tokenizer, we propose acoustic autoregressive modeling (AAR) which models the audio tokens with a next-scale paradigm. Since each scale contains multiple audio tokens, the AAR can lead to much fewer autoregressive step numbers during inference compared to the traditional token-level modeling. By reducing both the token length and the autoregressive step number, our approach achieves not only a superior generated audio quality but also a remarkably faster (about  $35 \times$ ) inference speed.

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

Our contribution is three-fold:

- We present Scale-level Audio Tokenizer (SAT) for audio reconstruction which can efficiently compress audio sequence to tokenizers with different scales.
- Based on SAT, we introduce scale-level Acoustic AutoRegressive modeling (AAR), significantly reducing the inference latency and training cost.
- · Extensive experiments are conducted to analyze the performance of the proposed approach, providing insights into its capabilities and potential applications in the field of audio synthesis.

#### 2 **Related Works**

**Raw audio discretization.** Before the development of Variational Autoencoders (VAEs) (Van Den Oord et al., 2017; Razavi et al., 2019), converting continuous domains into discrete representations was a significant challenge in generative modeling. VAEs facilitate the effective quantization of inputs into structured priors using powerful encoder-decoder networks, allowing manipulation in tasks like generation and understanding (Achiam et al., 2023; Touvron et al., 2023; Caillon and Esling, 2021). Recent innovations, such as VQGAN (Esser et al., 2021) and RQGAN (Lee et al., 2022), have further advanced these priors, improving model generalization and inspiring numerous works in audio discretization (Oord et al., 2016; Caillon and Esling, 2021; Siuzdak et al., 2024; Li et al., 2024b). In the audio domain, Encodec (Défossez et al., 2022) employs an architecture similar to SoundStream (Zeghidour et al., 2021), using an encoder-decoder model to reconstruct audio, incorporating residual quantization and a spectrogram-style discriminator to enhance audio quality. In contrast, HIFI-codec (Yang et al., 2023a) uses group residual quantization to refine the representation in the initial quantization layer. Kumar et al. (Kumar et al., 2024) have made significant contributions to audio reconstruction by



Figure 2: Our model involves two distinct training phases. **Stage 1**: Scale-level Audio Tokenizer (SAT) to encode an audio sample into a series of K tokens scales, donated as  $\mathcal{R} = (r_1, r_2, \ldots, r_K)$ . Each scale encodes information in different frequencies of the audio waveform. **Stage 2**: Acoustic AutoRegressive (AAR) modeling via next-scale prediction relies on the pre-trained SAT to predict each scale-level token  $r_i$  by conditioning on all previously predicted scales  $r_{<i}$  and a CLAP token (Wu et al., 2023b) as the start token. The CLAP token is derived from ground truth audio. During training, we use the standard cross-entropy loss and the attention mask as figured above to ensure that each  $r_i$  can only be attributed by  $r_{<i}$  and the start token.

introducing multi-spectrogram loss and quantizer dropout to enhance bitrate efficiency. Building upon these advances, our work poses an important question: can we use fewer tokens to represent lowfrequency information, thereby efficiently reducing the burden while maintaining high-quality audio reconstruction? To address this, we propose a Scalelevel Audio Tokenizer, which encodes audio on different scales, capturing hierarchical features that improve both the efficiency and quality of audio generation and reconstruction.

159

160

161

163

164

165

169

170

171

172

173

174

175

176

178

179

181

184

185

186

187

Autoregressive modeling The autoregressive model (Chowdhery et al., 2023; Hoffmann et al., 2022), as a different approach from diffusion models, leverages efficient Large Language Models (LLMs) (Vaswani et al., 2017; Devlin et al., 2018; Touvron et al., 2023; Achiam et al., 2023; Kreuk et al., 2022; Wu et al., 2023a) to generate the next tokens sequentially to construct the output. Due to its sequential nature, autoregressive models have excelled in text generation, machine translation, and other sequence prediction tasks. Recently, autoregressive models have also made significant processes in the image generation domain (Chang et al., 2022; Sun et al., 2024). By treating image pixels or patches as sequences, these models can generate high-quality images by sequentially predicting each part of the image.

## 3 Preliminary: Vanilla Audio Tokenizer

188Audio quantization. Consider an audio signal189 $a \in \mathbb{R}^{C \times T}$ , where C represents the number of190audio channels and T is the number of samples over191the duration of the signal. Traditional approach192(Kumar et al., 2024; Défossez et al., 2022; Yang

et al., 2023a) in audio tokenizer often involves a 1D convolutional-based autoencoder frameworks to compress audio waveform to latent space  $x \in \mathbb{R}^{l \times d}$  where *l* is the token length and then utilizes a vector quantization to quantize the latent tokens:

$$x = \mathcal{E}(a), \quad \hat{x} = \mathcal{Q}(x), \quad \hat{a} = \mathcal{D}(\hat{x})$$
 (1)

193

194

195

196

197

199

200

202

203

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

where  $\mathcal{E}(\cdot)$  donates encoder,  $\mathcal{Q}(\cdot)$  a vector quantizatier, and  $\mathcal{D}(\cdot)$  a decoder. A vector quantizer  $\mathcal{Q}$ maps each feature vector in the latent space x to the closest vector in a learnable codebook  $Z \in \mathbb{R}^{d \times V}$ with V vectors of dimension d. Specifically, vector quantization  $\hat{x} = \mathcal{Q}(x)$  involves looking up the closest match for each feature vector in x with vectors in Z by minimizing Euclidean distance:

$$\hat{x} = \operatorname{argmin}_{z \in Z} ||x - z||_2 \tag{2}$$

where  $\hat{x}$  represents the quantized output and x is the input to the quantizer.

However, due to the complexity of the audio waveform, particularly in handling frequencyspecific information, a residual quantization approach is typically employed. In residual quantization, a sequence of vector quantizers  $Q = \{Q_1, Q_2, \dots, Q_K\}$  is used, where each quantizer  $Q_i$  iteratively quantizes the residual error from the previous step. Specifically, after each quantization step, the residual error is computed as  $\delta_i = x_i - \hat{x}_i$  and passed to the next quantizer as the input  $x_i = \delta_{i-1}$ . The final quantized representation  $\hat{f}$  is obtained by summing the outputs from all quantizers

$$\hat{x} = \sum_{i=1}^{r} \hat{x_i} \tag{3}$$

which is then decoded by the decoder  $\mathcal{D}(\hat{x})$  to produce the reconstructed output  $\hat{a}$ .

**Loss function.** To train audio quantized autoencoder, we leverage a combination of loss functions including the reconstructed time-domain loss  $\mathcal{L}_t$ , reconstructed frequency domain loss  $\mathcal{L}_f$ , discriminative loss  $\mathcal{L}_G$ , residual quantization loss  $\mathcal{L}_{vq}$ , and commitment loss  $\mathcal{L}_{commit}$  (Défossez et al., 2022) :

$$\mathcal{L} = \lambda_t \mathcal{L}_t + \lambda_f \mathcal{L}_f + \lambda_G \mathcal{L}_G + \mathcal{L}_{vq} + \lambda_{com} \mathcal{L}_{com}.$$
(4)

Specifically, reconstructed time-domain loss measures the absolute difference between a and  $\hat{a}$  as

$$\mathcal{L}_t = ||a - \hat{a}|| \tag{5}$$

and frequency domain loss assesses the difference over mel-spectrograms across n time scales as

$$\mathcal{L}_{f} = \sum_{i=1}^{n} ||\mathcal{S}_{i}(a) - \mathcal{S}_{i}(\hat{a})|| + ||\mathcal{S}_{i}(a) - \mathcal{S}_{i}(\hat{a})||_{2}^{2}$$
(6)

where  $S_i$  represents the transformation to the melspectrogram at scale *i*. The discriminative loss is derived from a multi-scale STFT discriminator, as introduced in (Zeghidour et al., 2021) to ensure the model captures high-fidelity audio features across various time-frequency scales. The vector quantization loss encourages the encoded features to match the codebook vectors, and the commitment loss penalizes deviations from these vectors, ensuring that the encoder commits to the quantized space as

$$\mathcal{L}_{vq} = \sum_{i=1}^{r} ||\mathbf{sg}(x_i) - z_i||_2^2$$

$$\mathcal{L}_{com} = \sum_{i=1}^{r} ||x_i - \mathbf{sg}(z_i)||_2^2.$$
(7)

Analysis. The baseline audio tokenizer can successfully discretize audio tokens. However, due to the residual quantization, the token length representing each audio will be significant which severely hinders the efficiency in the autoregres-254 sive modeling. Considering each quantizer in resid-255 ual quantization basically divides and represents the audio into different frequency bands (Défossez et al., 2022; Kumar et al., 2024; Yang et al., 2023a), we aim to further adjust the token length based on its represented frequencies, i.e., lower-frequency parts can be represented with fewer tokens. To this end, we introduce the Scale-level Audio Tokenizer to reduce the number of tokens being used. 263

## 4 Method

Our approach consists of two major stages: (1) In the first stage, we train a Scale-level Audio Tokenizer (SAT) to convert continuous audio signals into discrete tokens using multi-scale residual quantization. (2) The second stage reformulates the Acoustic AutoRegressive modeling (AAR) in a next-scale manner and models the tokens obtained from the frozen SAT tokenizer with a transformer structure.

#### 4.1 Scale-level Audio Tokenizer

In Scale-level Audio Tokenizer (SAT), we employ the same encoder-decoder architecture as baseline tokenizer (Défossez et al., 2022) but incorporate multi-scale residual quantization (MSRQ) to enhance efficiency and flexibility in audio representation. In MSRQ, as shown in Fig. 2 (a), the quantizer  $Q_i$  is defined the same as the baseline setting as  $r_i = Q_i(r_{i-1})$  while the feature map  $r_{i-1}$  is first downsampled from its original dimension  $l_K \times d$ to a lower resolution  $l_k \times d$  where K is the scale number of the last index and k is the scale number of the correct index. After downsampling, the lookup procedure is performed to match each feature vector with the closest codebook vector  $Z_i$ . After the look-up, the processed quantized vector  $z_i$  is upsampled back to the original dimension  $l_K \times d$ to ensure consistency across scales. Due to the loss of high-frequency information from downsampling, we employ a 1D convolutional layer after upsampling to restore the lost details and enhance the fidelity of the reconstructed audio. Specifically, this convolutional layer processes the upsampled feature vectors according to the equation

$$\phi(\hat{r}) = \gamma \times \operatorname{conv}(\hat{r}) + (1 - \gamma) \times \hat{r} \qquad (8)$$

where  $conv(\cdot)$  applies a 1D convolution with a kernel size of 9. This design effectively combines the original features with the transformed outputs, while preserving the reparameterization inherent to vector quantization, controlled by the quantization residual ratio  $\gamma$ . In the Appendix, we provide a pseudo-code for the scale-level audio tokenizer.

#### 4.2 Acoustic AutoRegressive Modeling

**Vanilla autoregressive modeling.** Autoregressive modeling is first introduced by (Sutskever et al., 2014; Bahdanau et al., 2014) and quickly spread to different modalities such as image (Sun et al., 2024), video (Weissenborn et al., 2019) and

236

240

241

242

243

244

245

246

247

264 265

266

267

268 269

270 271 272

273

274

275

276

277

278

279

281

282

283

284

287

288

291

292

293

294

295

296

297

298

299

300

301

302

303

304

306

307

308

309

310

358 359

361

362

363

364

365

366

367

369

370

371

373

374

375

376

377

378

379

380

381

382

383

384

387

389

391

392

393

394

395

396

397

398

400

360

312 3D modeling (Siddiqui et al., 2024). In autoregres-313 sive modeling, a sequence of data points is modeled 314 as a product of conditional probabilities. For a se-315 quence  $x = (x_1, x_2, ..., x_T)$ , its joint distribution 316 can be expressed and modeled as

317

339

341

345

347

353

354

357

$$p(x_1, x_2, ..., x_T) = \prod_{t=1}^T p(x_t | x_1, x_2, .... x_{t-1}).$$
(9)

This approach is widely used across various domains due to its flexibility and ability to capture dependencies within data. For any continuous modal-321 ity, it is traditional to first train a tokenizer to discretize the input into tokens, which can then be modeled using a discrete categorical distribution. 323 This step involves mapping the continuous data to a sequence of discrete tokens  $x = (x_1, x_2, \ldots, x_T)$ that are fed into an autoregressive model to pre-326 dict the next token in the sequence, based on the preceding tokens. In the context of transformers, which have become the dominant architecture for autoregressive modeling, the attention mechanism plays a crucial role in training. The attention mechanism allows the model to focus on different parts of the input sequence when making predictions. To 333 ensure that the model adheres to the autoregressive 334 property, where each token  $x_i$  is predicted based only on previous tokens  $x_1, x_2, \ldots, x_{i-1}$ , an attention mask is applied. Mathematical, the attention 338 mask M is defined as

$$M_{ij} = \begin{cases} 1, & \text{if } i \le j \\ 0, & \text{otherwise} \end{cases}$$
(10)

where guarantee the modeling's performance on predicting  $x_i$  is only relevant to its preceding tokens.

After the completion of training of such a model P using cross-entropy loss, it can efficiently handle complex dependencies and generate new samples by sequentially predicting each token conditioned on its predecessors (Achiam et al., 2023; Touvron et al., 2023; Sun et al., 2024).

This capability makes autoregressive models well-suited for generating data that requires a coherent and consistent sequence. However, their capacity for audio generation is still under-explored due to the huge sequence length required for audio data. The sheer number of tokens needed to represent even short audio clips can lead to computational inefficiencies and challenges in maintaining temporal coherence. To efficiently solve such a challenge, we combine the unique property of our SAT to efficiently generate audio via scale-level Acoustic AutoRegressive modeling.

Acoustic autoregressive modeling. To shorten the inference step, we propose Acoustic autoregressive modeling (AAR). This approach, distinct from traditional vanilla autoregressive models that predict token sequences one by one, involves predicting across different scales. Attributed by SAT, our method represents an audio sample as a series of scale-level representations:

$$\mathcal{R} = (r_1, r_2, \dots, r_K) \tag{11}$$

By efficiently expressing it as joint modeling, the audio sequence is defined as:

$$p(\mathcal{R}) = \prod_{i=1}^{K} p(r_i | r_1, r_2, ..., r_{i-1})$$
(12) 372

In this formulation, each  $r_i$  represents a distinct scale in the hierarchical representation of the audio signal. The model predicts each scale by conditioning on all previously predicted scales, effectively capturing both global structures and fine-grained details of the audio. This hierarchical approach reduces the complexity associated with long sequence lengths by leveraging multi-scale dependencies, thereby enhancing the model's efficiency and ability to maintain temporal coherence. To successfully implement our method, we modify the attention mask M for each scale  $r_i$  to focus only on the relevant scales:

$$M_{ij} = \begin{cases} 1, & \text{if } i \le j \\ 0, & \text{otherwise} \end{cases}$$
(13)

This attention mask ensures that the model only attends to  $r_1, r_2, \ldots r_{i-1}$  when predicting  $r_i$  ignoring future scales and reducing unnecessary computations. A detailed description of the implementation is summarized in the Appendix.

## 5 Experiment

## 5.1 Evaluation Metrics and Settings

We evaluate FAD (Kilgour et al., 2018), MEL distance (Kumar et al., 2024), and STFT distance (Kumar et al., 2024) as reference for reconstruction, and FAD (Kilgour et al., 2018), ISc and KL (Salimans et al., 2016) for generation. FAD, built upon VGGish (Chen et al., 2020), is the metric to indicate the similarity of the generated and target

ID	Mathad	Dataset	Reconstruction			Generation	
ID	Method		$ \mathcal{L} $	rFAD↓	MEL↓	gFAD↓	KL↓
1	Vanilla AR (baseline)	AS	750	1.39	1.33	10.05	3.01
2	DiffSound (Yang et al., 2023b)	AS	-	-	-	9.76	4.21
3	AudioLCM (Liu et al., 2024b)	AC + LP	-	-	-	3.92	1.20
4	AudioLDM 2 (Liu et al., 2024a)	AS + 7 more	-	-	-	1.42	0.98
5	AAR (ours)	AS	455	1.09	1.33	6.01	2.27

Table 1: We evaluated the performance of our AAR model against other methods using rFAD and MEL Distance to measure reconstruction quality, and gFAD and KL Divergence to assess generation quality, where  $\downarrow$  indicates that lower values are better. In this context,  $|\mathcal{L}|$  denotes the token length. Additionally, AS, AC, and LP denote the datasets AudioSet, AudioCaps, and LP-Musicaps, respectively.



Figure 3: Performance of autoregressive model when classifier-free guidance is 10. next-token: AR via next-token prediction; next-scale: our AAR.

samples effectively. MEL distance quantifies the difference in mel-spectrogram features, and STFT distance measures the short-time Fourier transform discrepancies between the generated and target audio signals, which focus more on high-frequency information for audio. Additionally, ISc, simulating its performance on image generation, is used to evaluate the generated sample diversity and quality. KL divergence is utilized to measure the difference between the probability distributions of the generated and target samples.

401

402

403

404

405

406

407

408

409

410

411

We conducted all experiments on the AudioSet 412 (Gemmeke et al., 2017) dataset. To effectively eval-413 uate the performance of our audio tokenizer, we 414 divided the original 10-second evaluation set into 415 n segments, each matching the window size of 416 417 our model for reconstruction. After reconstructing these segments, we reassembled them into a 418 complete audio stream. For autoregressive genera-419 tion, we randomly selected one segment from the 420 evaluation set and used it as the ground truth. 421

#### 5.2 Implementation Details

**Tokenizer.** In stage 1, we utilize multi-scale residual quantization (MSRQ) of codebook size 1024 with the Soundstream autoencoder framework (Zeghidour et al., 2021). The model is trained for 100 epochs using the Adam optimizer with  $\beta_1 = 0.5$  and  $\beta_2 = 0.9$ . We apply a cosine learning rate scheduler with initial learning rate 3e-4 and set the loss weights to  $\lambda_t = 0.1$ ,  $\lambda_f = 3$ ,  $\lambda_G = 3$ , and  $\lambda_{com} = 1$ . Our discriminator updated 2/3 times during training.

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

**Transformer.** In stage 2, our primary focus is on scale-level acoustic autoregressive modeling. To achieve this, we employ a GPT-2-style transformer (Radford et al., 2019) with adaptive normalization (Zhang et al., 2018) and depth of 16. We utilize CLAP audio embeddings (Wu et al., 2023b) as our start tokens. Since one-second audio segments often contain limited meaningful information, we opt to use 10-second audio embeddings to capture richer context, even when generating one-second clips. For training, we adopt the AdamW optimizer with a learning rate of 1e-4, using a linear learning rate scheduler. Additionally, we apply a weight decay of 0.05 and implement warmup settings with an initial warmup proportion of 0.005 and an end warmup proportion of 0.01.

#### 5.3 Main Results Analysis

We compare the performance of our approach with previous approaches. As shown in Tab 1, our proposed SAT tokenizer suppresses the baseline encodec (Défossez et al., 2022) by 0.3 FAD in the reconstruction task, despite using fewer tokens (750 tokens v.s. 455 tokens). This shows that by increasing quantization while reducing the number of tokens, we can efficiently improve reconstruction quality while using fewer tokens.

# scales	# tokens	FAD	MEL	STFT
10	207	1.81	1.55	2.08
16	303	1.52	1.46	1.80

Table 2: Ablation study of SAT performance in number of scales.

Scheduler	# tokens	FAD	MEL	STFT
Logarithmic	303	1.52	1.46	1.80
Quadratic	455	1.40	1.37	1.86
Linear	601	1.38	1.39	2.02

Table 3: Ablation study on scale setting of SAT. All scale settings are trained in the same numbers of scale.

For the audio generation, we introduce an autoregressive model with next-token prediction as the baseline. To ensure a fair comparison, we employ two encoders (Encodec (Défossez et al., 2022) for AR and our SAT for AAR) with similar performance. We find that our proposed AAR shows superior performance in terms of both latency and audio quality. As shown in Fig. 3, the next-scale prediction demonstrates a remarkable improvement (0.225s v.s. 7.866s) and generation enhancement (FAD 5.55 v.s. 6.88). More analysis of training costs is available in the Appendix.

#### 5.4 Ablation Experiments

We conduct ablation experiments to validate the effectiveness of the components in SAT and AAR.

**Effect of discriminator** We explored multiple discriminator configurations to optimize the performance of our Scale-level Audio Tokenizer (SAT). As illustrated in Tab 6, we tested two different discriminator setups: one using only a multi-scale short-time Fourier transform (STFT) discriminator (Zeghidour et al., 2021) and another combining the multi-scale STFT discriminator with a multi-

Latent dim.	FAD	MEL	STFT
8	1.47	1.55	2.15
16	1.38	1.52	2.14
32	1.60	1.43	2.05
64	1.09	1.33	1.98

Table 4: Ablation study on latent dimension. We fix the scale to 16 and use the same quadratic scale setting. "Latent dim." represents dimension of latent representation.

Window	FAD	MEL	STFT
1s	1.22	1.36	1.85
5s	1.29	1.41	1.93

Table 5: Ablation study on temporal window.

STFTD	MPD	MSD	FAD	MEL	STFT
√	X	X	1.38	1.36	1.76
$\checkmark$	$\checkmark$	$\checkmark$	2.29	1.65	2.12

Table 6: Ablation study on discriminator choice. STFTD stands for Multi-scale short-time fourier transform discriminator, MPD stands for multi-period discriminator, MSD stands for multi-scale discriminator.

period discriminator (MPD) (Kong et al., 2020) and a multi-scale discriminator (MSD) (Kumar et al., 2019). Our results indicate that using only a multiscale STFT discriminator is sufficient for effective reconstruction. 483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

Effect of the scale setting. To find the optimal combination of SAT configuration, we start with Encodec in 128 latent dimensions with 10 quantizers (Défossez et al., 2022) and test multiple scales with shared codebooks of different sizes and individual codebooks for each scale. In particular, Tab 2 shows that enlarging the scale to 16 consistently improved audio quality. As illustrated in Tab 3, we tested the performance of linear, quadratic and logarithmic scheduling on 16 scales: linear scheduling provides a balanced number of tokens for each scale; quadratic scheduling focuses more on the early or late stages of the process; and logarithmic scheduling offers a more gradual progression. We believe the suboptimal performance observed in logarithmic scheduling is due to its lack of high-frequency information representation at larger scales even though it also builds a complete information flow for audio. Quadratic scheduling, in particular, proved to be more efficient, requiring fewer tokens than linear scheduling (455 v.s. 601) and also achieves comparable reconstruction performance in audio quality.

To further improve the model's capacity, we fixed the decoder dimension to 1024 and tested latent dimensions of 8, 16, 32, and 64. As Tab 4 indicated, our SAT achieves its superior performance in the latent dimension of 64.

**Effect of temporal windows change.** To effectively validate the performance of our scale schedul-

478

479

480

481

482

ID	Description	FAD↓	IS↑	KL↓	Latency $\downarrow$ (s)
1	Vanila AR	10.05	2.42	3.01	7.86
2	AAR	$9.24_{-0.81}$	$2.69_{+0.27}$	$2.94_{-0.07}$	$0.21_{-7.21}$
3	+ Attn. Norm	$8.80_{-1.25}$	$2.80_{\pm 0.38}$	$2.79_{-0.22}$	$0.25_{-7.61}$
4	+ CFG	$6.44_{-3.61}$	$3.52_{\pm 0.90}$	$2.32_{-0.69}$	$0.25_{-7.61}$
5	+ Top-k	$6.25_{-3.81}$	$3.59_{+1.17}$	$2.30_{-0.71}$	$0.25_{-7.61}$
6	+ Top-p	$6.01_{-4.04}$	$3.68_{\pm 1.26}$	$2.27_{-0.74}$	$0.25_{-7.61}$

Table 7: Ablation study on components of AAR. vanilla AR and AAR are implemented in GPT-2 style transformer with adaptive layer normalization; "Attn. Norm" represents normalizing q and k into unit vector before attention; "CFG" means classifier free guidance scale of 2; Top-k and Top-p are sampling strategies where Top-k randomly selects from the top 200 indices, and Top-p (nucleus sampling) selects tokens with a cumulative probability of 0.95.



Figure 4: Performance of AAR in different classifierfree guidance scales from 2 to 18 (left to right), with each point incremented by 2. The red line represents Fréchet Audio Distance (FAD) v.s. Inception Score (ISc), while the blue line represents Kullback-Leibler divergence (KL) vs. Inception Score (ISc).

ing in audio reconstruction, we train our SAT using 518 5-second audio windows with the  $5 \times$  original quan-519 tizer setting. This approach allows us to assess our 520 SAT's ability to handle varying temporal dimen-521 sions and capture essential audio features over dif-522 ferent time scales. By experimenting with different 523 window sizes, we aim to determine the optimal 524 configuration for maintaining high reconstruction 525 quality while maximizing efficiency. The results of these experiments are presented in Tab 5. we find that the reconstruction quality between 1-second 528 and 5-second windows is similar, suggesting that our SAT performs well across diverse time windows, maintaining consistent quality and demon-531 strating robustness in handling varying temporal scales. 533

Effect of AAR and sampling technique. We
evaluate our AAR with the same setting as the
baseline vanilla AR with roadmap shown in Tab 7.
We notice that our AAR can not only improve the
generation abilities but also significantly reduce

the inference time to an acceptable range. Moreover, the introduction of attention normalization can stabilize the training and further enhance the model's performance, leading to improved FAD and IS scores. The addition of CFG and advanced sampling techniques such as top-k and top-p sampling continues to push the boundaries of audio generation quality. 539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

**Effect of classifier-free guidance.** As shown in Fig. 4, we evaluate the relationship between the Inception Score with Fréchet Audio Distance and Inception Score with KL divergence across different Classifier-Free Guidance scales. We find that as the CFG scales increase, the ISc improves, while both FAD and KL metrics converge and stabilize at CFG = 14 and finally achieving FAD 5.19.

## 6 Conclusion

In this paper, we introduced a novel approach for audio generation using a multi-scale autoregressive model via next-scale prediction. This framework leverages the scale-level audio tokenizer, which efficiently compresses audio sequences into tokenizers of varying scales, thereby improving efficiency while maintaining high fidelity. Through comprehensive experiments, we demonstrated the superior performance of our method in generating high-quality audio compared to traditional autoregressive methods.

Our approach provides an efficient solution for audio generation. By incorporating a multi-scale residual quantization technique, the model effectively reduces the sequence length required for generation, leading to enhanced efficiency and reduced computational demands.

## 7 Limitations.

573

593

594

595

596

598

603

605

608

610

611

612

613

614

615

616

617

618

619

621

623

Despite the strong performance of next-scale prediction in general audio generation with text con-575 trol, several limitations remain that warrant further 576 exploration. Signal-level audio tokenizers, such as those employed in this work, often rely on residual quantization to capture information across different frequencies. While effective, this approach faces challenges in managing long token lengths, particularly for high-resolution audio signals. To mitigate this, we adopt a multi-scale approach that reduces token length while maintaining reconstruction qual-584 ity. However, semantic tokenizers offer a promis-585 ing alternative by achieving shorter token lengths with higher information density. Integrating semantic information into multi-scale quantized tokens 588 could further reduce token length while enhancing the richness and efficiency of latent representations. Addressing this integration and improving scalability will be a key direction for our future research.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
  - Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Antoine Caillon and Philippe Esling. 2021. Rave: A variational autoencoder for fast and highquality neural audio synthesis. *arXiv preprint arXiv:2111.05011*.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. 2022. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325.
- Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. 2020. Vggsound: A large-scale audio-visual dataset. In ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 721–725. IEEE.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul

Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113. 624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

- Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. 2022. High fidelity neural audio compression. *arXiv preprint arXiv:2210.13438*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020.
  An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883.
- Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. 2017. Audio set: An ontology and human-labeled dataset for audio events. In 2017 IEEE international conference on acoustics, speech and signal processing (ICASSP), pages 776–780. IEEE.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*.
- Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM transactions on audio, speech, and language processing*, 29:3451–3460.
- Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. 2018. Fr\'echet audio distance: A metric for evaluating music enhancement algorithms. *arXiv preprint arXiv:1812.08466*.
- Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. *Advances in neural information processing systems*, 33:17022– 17033.
- Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. 2022. Audiogen: Textually guided audio generation. *arXiv preprint arXiv:2209.15352*.

- Kundan Kumar, Rithesh Kumar, Thibault De Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre De Brebisson, Yoshua Bengio, and Aaron C Courville. 2019. Melgan: Generative adversarial networks for conditional waveform synthesis. *Advances in neural information processing systems*, 32.
- Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar. 2024. Highfidelity audio compression with improved rvqgan. *Advances in Neural Information Processing Systems*, 36.
- Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. 2022. Autoregressive image generation using residual quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11523–11532.

691

703

704

705

706

708

710

711

712

713

715

716

719

721

726

727

729

731

732

- Xiang Li, Kai Qiu, Hao Chen, Jason Kuen, Jiuxiang Gu, Bhiksha Raj, and Zhe Lin. 2024a. Imagefolder: Autoregressive image generation with folded tokens. *arXiv preprint arXiv:2410.01756*.
- Xiang Li, Kai Qiu, Hao Chen, Jason Kuen, Jiuxiang Gu, Jindong Wang, Zhe Lin, and Bhiksha Raj. 2024b. Xq-gan: An open-source image tokenization framework for autoregressive generation. *arXiv preprint arXiv:2412.01762*.
- Xiang Li, Kai Qiu, Hao Chen, Jason Kuen, Zhe Lin, Rita Singh, and Bhiksha Raj. 2024c. Controlvar: Exploring controllable visual autoregressive modeling. *arXiv preprint arXiv:2406.09750*.
- Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley. 2024a. Audioldm 2: Learning holistic audio generation with selfsupervised pretraining. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- Huadai Liu, Rongjie Huang, Yang Liu, Hengyuan Cao, Jialei Wang, Xize Cheng, Siqi Zheng, and Zhou Zhao. 2024b. Audiolcm: Text-to-audio generation with latent consistency models. *arXiv preprint arXiv:2406.00356*.
- Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. 2016. Wavenet: A generative model for raw audio. arXiv preprint arXiv:1609.03499.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. 2019. Generating diverse high-fidelity images with vq-vae-2. Advances in neural information processing systems, 32.

Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. 2016. Improved techniques for training gans. *Advances in neural information processing systems*, 29.

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

765

767

768

769

770

774

775

776

777

779

780

781

782

783

784

785

786

- Yawar Siddiqui, Antonio Alliegro, Alexey Artemov, Tatiana Tommasi, Daniele Sirigatti, Vladislav Rosov, Angela Dai, and Matthias Nießner. 2024. Meshgpt: Generating triangle meshes with decoder-only transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19615–19625.
- Hubert Siuzdak, Florian Grötschla, and Luca A Lanzendörfer. 2024. Snac: Multi-scale neural audio codec. *arXiv preprint arXiv:2410.14411*.
- Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. 2024. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.
- Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. 2024. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *arXiv preprint arXiv:2404.02905*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural discrete representation learning. *Advances in neural information processing systems*, 30.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. 2019. Scaling autoregressive video models. *arXiv preprint arXiv:1906.02634*.
- Tong Wu, Zhihao Fan, Xiao Liu, Hai-Tao Zheng, Yeyun Gong, Jian Jiao, Juntao Li, Jian Guo, Nan Duan, Weizhu Chen, et al. 2023a. Ar-diffusion: Autoregressive diffusion model for text generation. *Advances in Neural Information Processing Systems*, 36:39957–39974.
- Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. 2023b. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In *ICASSP 2023-2023 IEEE International*

788 Conference on Acoustics, Speech and Signal Process 789 ing (ICASSP), pages 1–5. IEEE.

790

791 792

793

794

795

796

797

798

802

803 804

805

806

807

808 809

810

811

812

- Dongchao Yang, Songxiang Liu, Rongjie Huang, Jinchuan Tian, Chao Weng, and Yuexian Zou. 2023a. Hifi-codec: Group-residual vector quantization for high fidelity audio codec. *arXiv preprint arXiv:2305.02765*.
- Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. 2023b.
   Diffsound: Discrete diffusion model for text-tosound generation. *IEEE/ACM Transactions on Audio*, *Speech, and Language Processing*, 31:1720–1733.
- Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. 2021. Vectorquantized image modeling with improved vqgan. *arXiv preprint arXiv:2110.04627*.
  - Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. 2021. Soundstream: An end-to-end neural audio codec. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:495–507.
- Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. 2018. Self-attention generative adversarial networks. In *International Conference* on Machine Learning, pages 7354–7363.

Algorithm 1 Multi-scale residual quantization **Input**: Raw Audio Signal  $\mathcal{A}$ **Parameter**: Encoder  $\mathcal{E}$ , Decoder  $\mathcal{D}$ , Quantizer  $\mathcal{Q}_{i=1}^{K}$ Hyper-parameter: multi-scale Resolution  $(t_k)_{k=1}^K$ , interpolation method  $\phi$ 1:  $f = \mathcal{E}(\mathcal{A})$ 2: R = []3: for k in (1, 2, ..., K) do if k = K then 4:  $r_k = \mathcal{Q}_k(f)$ 5:  $z_k = lookup(Q_k, r_k)$ 6: 7: else  $r_k = \mathcal{Q}_k(interpolate(f, t_k))$ 8:  $z_k = lookup(Q_k, r_k)$ 9:  $z_k = interpolate(z_k, t_K)$ 10: end if 11:  $R = R + r_k$ 12:  $f = f - \phi(z_k)$ 13: 14: **end for** 15: return R

### A Supplementary material

814

816

818

820

821

823

827

831

833

836

838

#### A.1 Multi-Scale Residual Quantization

Multi-scale residual quantization (MSRQ) in our SAT is designed to efficiently encode audio signals by leveraging multiple quantization stages on different scales. Specifically, the MSRQ process begins with the raw audio signal being encoded into a feature representation. This representation is then passed through a series of quantizers, each corresponding to a different scale. For each scale, the feature map is downsampled to match the target resolution before quantization. After quantization, the feature map is upsampled back to its original resolution. Due to information loss in interpolation, the upsampled feature map is further processed through our upsampling network to recover the information for each scale. The residual error, calculated after each quantization step, is passed to the next quantizer, allowing the model to refine the audio representation iteratively. The pesudo-code implementation can be shown in Algorithm 1.

## A.2 Scale-level Acoustic AutoRegressive Generation.

Our AAR method begins by taking the scale-level tokens generated by the MSRQ process. These to-

Algorithm 2 Multi-scale AR Generation

# Input: Text $\mathcal{T}$

**Parameter**: Decoder  $\mathcal{D}$ , GPT AR, conditional Model C

**Hyper-parameter**: multi-scale Resolution  $(t_k)_{k=1}^K$ , interpolation method  $\phi$ 

1:  $x_0 = C(T)$ 2:  $R = [], S = [x_0]$ 3: for k in (1, 2, ..., K) do  $x_k = AR(S)$ 4:  $R = R + x_k$ 5: if k = K then 6: 7: break 8: else  $x_k = interpolate(x_k, t_K)$ 9:  $x_k = \phi(interpolate(x_k, t_{k+1}))$ 10: 11: end if  $S = S + x_k$ 12: 13: end for 14: A = D(R)15: return A

kens are organized hierarchically, with each scale capturing different levels of detail in the audio signal-from coarse, low-frequency information to fine, high-frequency details. In generation, the process is structured to predict these scales sequentially, starting from the coarsest scale and progressing to finer scales. As Algorithm 2 illustrated, our AAR first initializes the generation process by producing an initial latent representation from the input text using a conditional model. This initial representation serves as the starting point for the autoregressive prediction. The model then proceeds through each scale, beginning with the coarsest, and generates the corresponding tokens by conditioning on the sequence of tokens generated thus far. After each scale's tokens are predicted, they are interpolated and refined to align with the resolution of the next finer scale. This iterative process continues until all scales are generated, ensuring a smooth and coherent progression from low to highfrequency details. Finally, the aggregated tokens from all scales are decoded into a complete audio signal, resulting in a high-fidelity output that effectively captures the nuances of the original audio.

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

## A.3 Scale Scheduling

In our paper, we explore three types of scale scheduling: Linear, Quadratic, and Logarithmic.

Specifically, Linear scheduling ensures that the difference between each scale is consistent and linear. For example, in our linear scheduling approach, we start from a scale of 1 and increase to 75 using 16 scales, resulting in a difference of approximately 5 between each consecutive scale. The detail visualization of our scale setting can be shown in Fig. 5.

867

868

870

871

872

874

876

878

881

885

892

896

900



Figure 5: Visualization of Linear, Quadratic, and Logarithmic scale scheduling across the range from 1 to 75.

#### A.4 Codebook Utilization

We analyze the codebook utilization across different models. In particular, we observe that the codebook utilization for Encodec (Défossez et al., 2022) consistently reaches 99%, indicating that the entire codebook is actively used during the encoding process. In contrast, our SAT model exhibits a lower utilization rate. Specifically, we find that the codebook utilization in every scale of the SAT model remains at approximately 60%. We hypothesize that this discrepancy is caused by the inherent structure of the SAT model, where each time we downsample the input and apply quantization, the model becomes increasingly selective in its use of the codebook entries.

A.5 Training Cost Analysis

Our AAR exhibits strong performance in terms of latency and quality while also efficiently reducing the training cost of the model. To be more specific, the vanilla AR achieved an FAD of 8.07 with a classifier-free guidance scale of 4 after training for 100 epochs, while our AAR achieved an improved FAD of 5.70 under the same settings in just 45 epochs. As shown in Tab 8, to achieve the same capacity of vanilla AR, our AAR only needs to train 20 epochs and efficiently save approximately 80% training cost.

Epoch	Method	FAD↓	IS↑	KL↓
100	AR	7.19	2.78	2.73
10	AAR	$7.57_{\pm 0.38}$	$2.97_{\pm 0.19}$	$2.74_{+0.01}$
20	AAR	$6.83_{-0.36}$	$3.24_{\pm 0.46}$	$2.53_{-0.20}$
30	AAR	$6.36_{-0.83}$	$3.49_{+0.71}$	$2.40_{-0.33}$
40	AAR	$6.32_{-0.87}$	$3.55_{\pm 0.77}$	$2.32_{-0.41}$
45	AAR	$6.13_{-1.06}$	$3.63_{\pm 0.85}$	$2.28_{-0.45}$

Table 8: Comparison of training cost between vanilla AR and our AAR. All results are generated within classifier-free guidance scale of 4.



Figure 6: Ablation study of upsampling functions on SAT.

#### A.6 Upsampling Function

In our work, to efficiently recover information loss from downsampling, we use a 1D convolutional layer after vanilla upsampling to ensure unique information on each scale is preserved and accurately represented. We evaluated its effectiveness through three configurations: unshared (each quantizer has its own convolutional layer); partially shared (approximately three quantizers share one layer); and fully shared (all quantizers use the same layer) to validate the effectiveness of this approach in distinguishing and splitting information across different



Figure 7: Ablation study of upsampling functions on AAR.

903

910

911

913 scales for multi-scale reconstruction, and so on, generation. Our experiments (see Fig. 6) show that 914 the performance of unshared, partially shared, and 915 fully shared networks in reconstruction is similar, 916 indicating that all configurations effectively main-917 918 tain audio quality during reconstruction. However, their impact on generation can be seen in Fig. 7, 919 where the partially shared architecture significantly 920 improves generation quality. 921